

To get meaningful results, all proteins must be compared on an apples-to-apples basis. We do this by fixing a window size k (typically 50-500 amino acid residues) and comparing all proteins using only a k -mer from each protein. This would be very difficult to do efficiently with BLAST. Our algorithm takes a set of organisms and finds the most conserved (orthologous families of) proteins, ranking them from most to least conserved. The conservation measure for each family is the “diameter” of the cluster formed by selecting a k -mer from each protein in the family. The algorithm computes the k -mers that give the maximum amount of conservation.

[illegible][illegible]

Some applications of an algorithm for computing protein conservation are:

- *Phylogeny.* Phylogenetic trees are traditionally built with rRNA. Using highly conserved proteins instead yields higher quality trees (*submitted for publication*).
- *Finding characteristic proteins for a clade.* The algorithm can be run over specific clades, and the proteins which are significantly more highly conserved compared to a superclade can be identified.